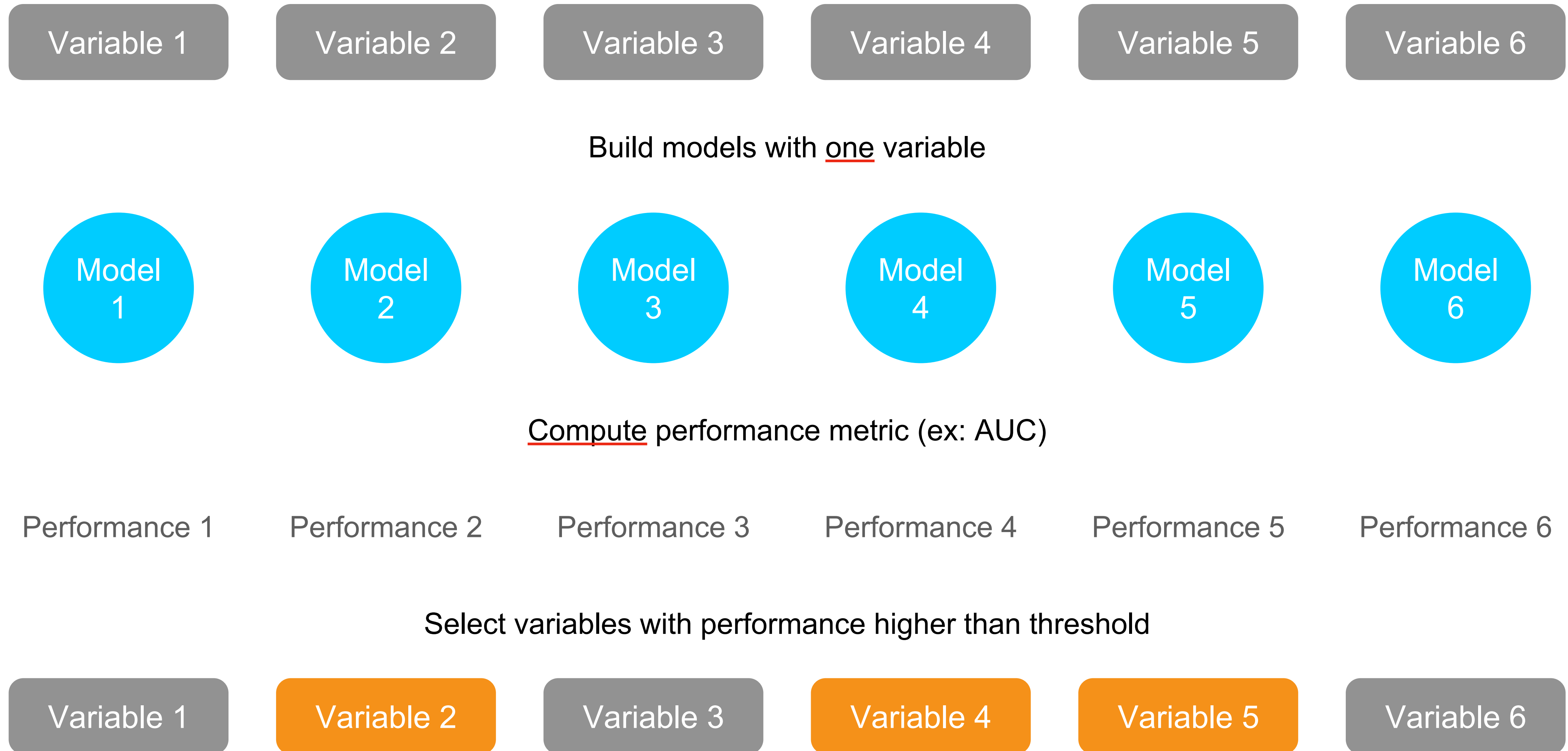Feature selection

# Feature selection

- Typical predictive analytical projects have base tables with ~ 1000 variables

- Which problems arise if we use all of these in our model?

    - Model very <u>hard</u> to interpret

    - Impossible to <u>present</u> to business

    - <u>Unstable</u> on long term (all variables need to be up to date)

    - Unnecessary <u>complexity</u> (takes longer to score model)

    - <u>Overfitting</u>

# Feature selection

- Decision trees: variable selection is <u>incorporated</u> in modeling
- Other algorithms
  - Univariate selection
  - Stepwise selection
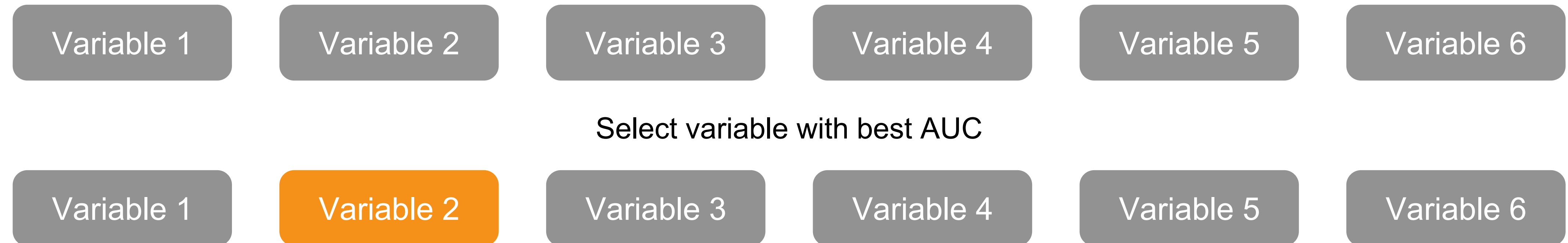
# Univariate variable selection

| Variable 1 | Variable 2 | Variable 3 | Variable 4 | Variable 5 | Variable 6 |

Build models with <u>one</u> variable

Model 1    Model 2    Model 3    Model 4    Model 5    Model 6

<u>Compute</u> performance metric (ex: AUC)

| Performance 1 | Performance 2 | Performance 3 | Performance 4 | Performance 5 | Performance 6 |

Select variables with performance higher than threshold

| Variable 1 | Variable 2 | Variable 3 | Variable 4 | Variable 5 | Variable 6 |

# Univariate variable selection

**Alternative methods**

- Information gain

- Pearson correlation with target higher than threshold

- Hypothesis testing

- ...

# Stepwise forward variable selection

## Step 1 : Model with 1 variable

| Variable 1 | Variable 2 | Variable 3 | Variable 4 | Variable 5 | Variable 6 |

Select variable with best AUC

| Variable 1 | Variable 2 | Variable 3 | Variable 4 | Variable 5 | Variable 6 |

# Stepwise forward variable selection

## Step 2 : Model with 2 variables

| | | | | |
|---|---|---|---|---|
| Variable 2<br>Variable 1 | | Variable 2<br>Variable 3 | Variable 2<br>Variable 4 | Variable 2<br>Variable 5 | Variable 2<br>Variable 6 |

Select variables with best AUC

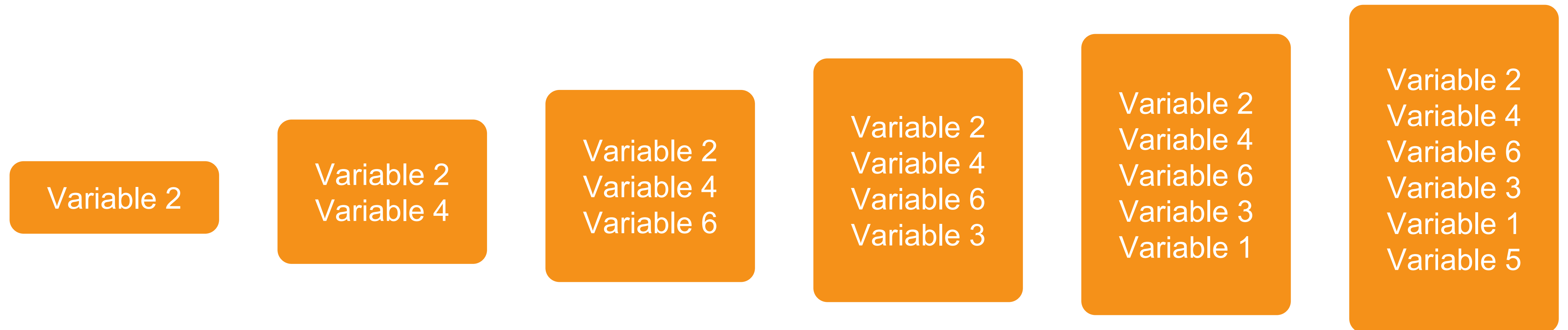| | | | | |
|---|---|---|---|---|
| Variable 2<br>Variable 1 | | Variable 2<br>Variable 3 | Variable 2<br>Variable 4 | Variable 2<br>Variable 5 | Variable 2<br>Variable 6 |

# Stepwise forward variable selection

## Step 3 : Model with 3 variables



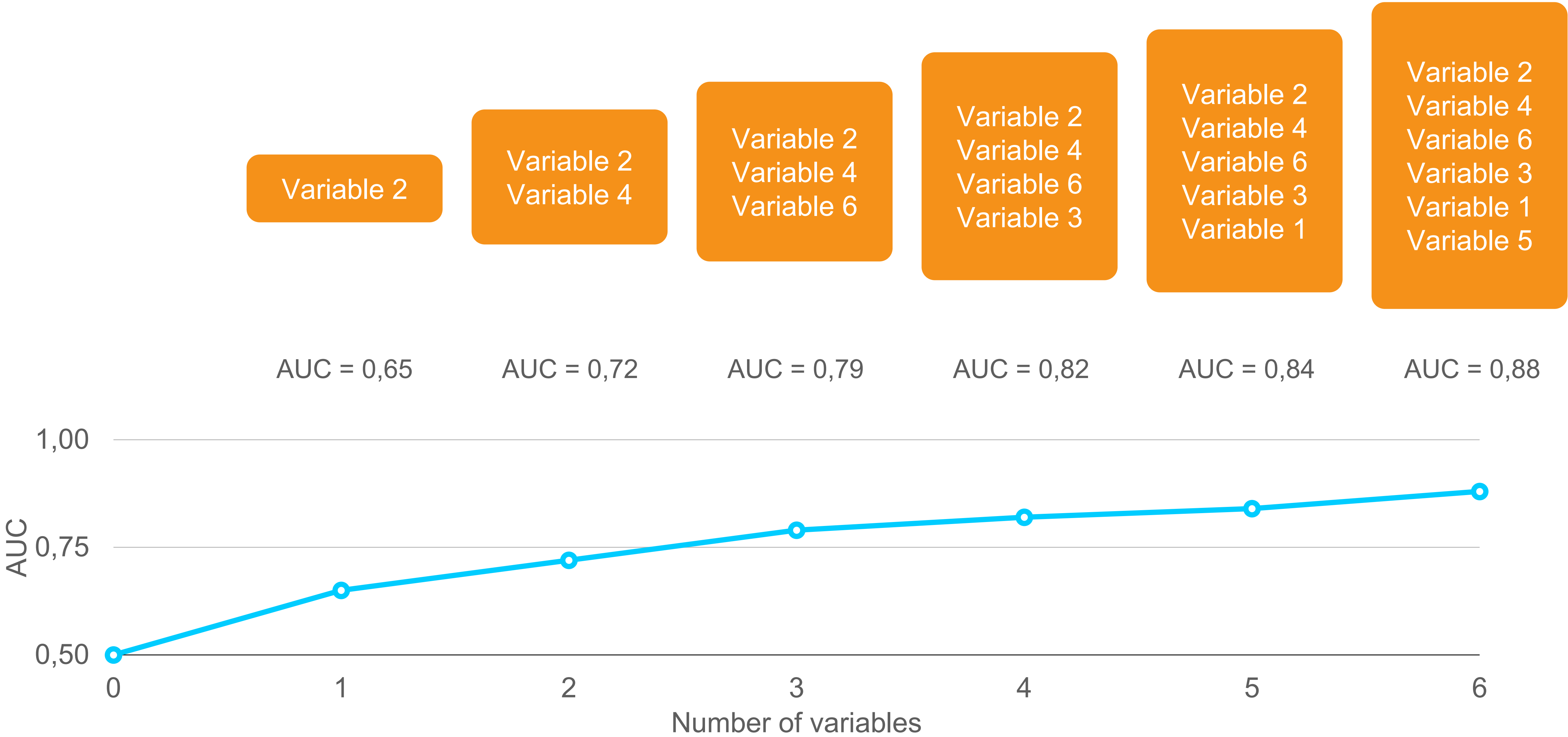| | | | |
|---|---|---|---|
| Variable 2<br>Variable 4<br>Variable 1 | Variable 2<br>Variable 4<br>Variable 3 | Variable 2<br>Variable 4<br>Variable 5 | Variable 2<br>Variable 4<br>Variable 6 |

Select variables with best AUC

| | | | |
|---|---|---|---|
| Variable 2<br>Variable 4<br>Variable 1 | Variable 2<br>Variable 4<br>Variable 3 | Variable 2<br>Variable 4<br>Variable 5 | Variable 2<br>Variable 4<br>Variable 6 |

# Stepwise forward variable selection
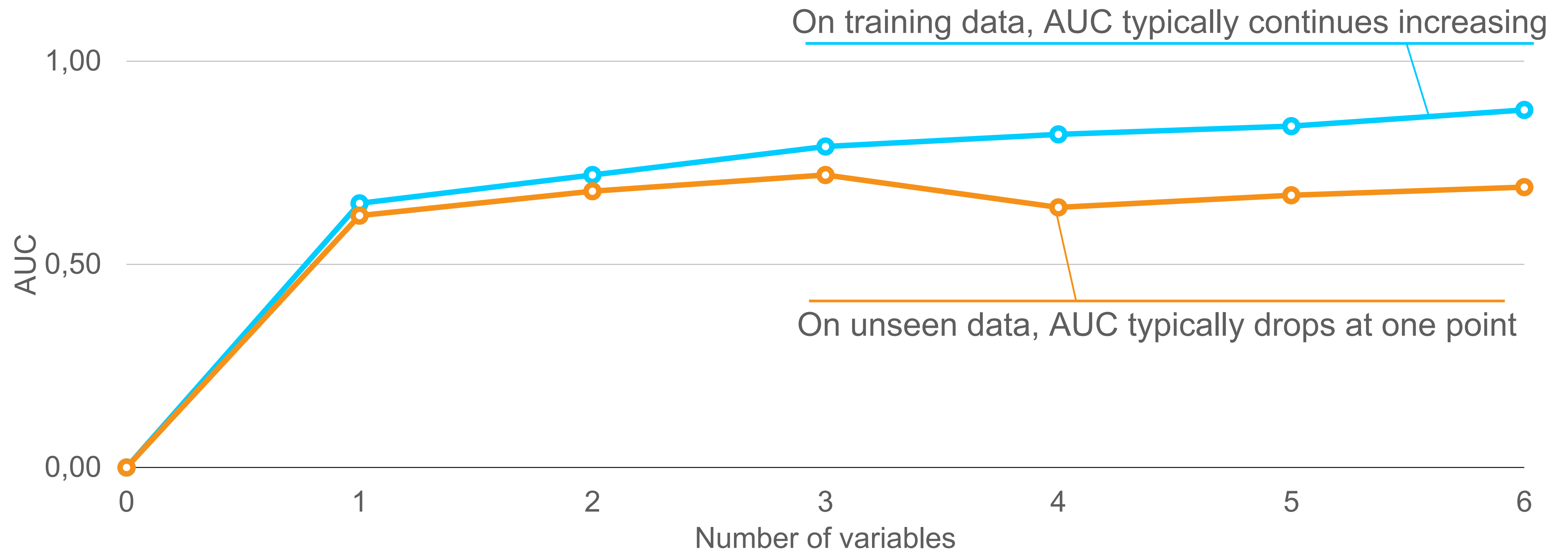
## Step N : Model with all variables

Result : N models, each with 1 additional variable

| Variable 2 |

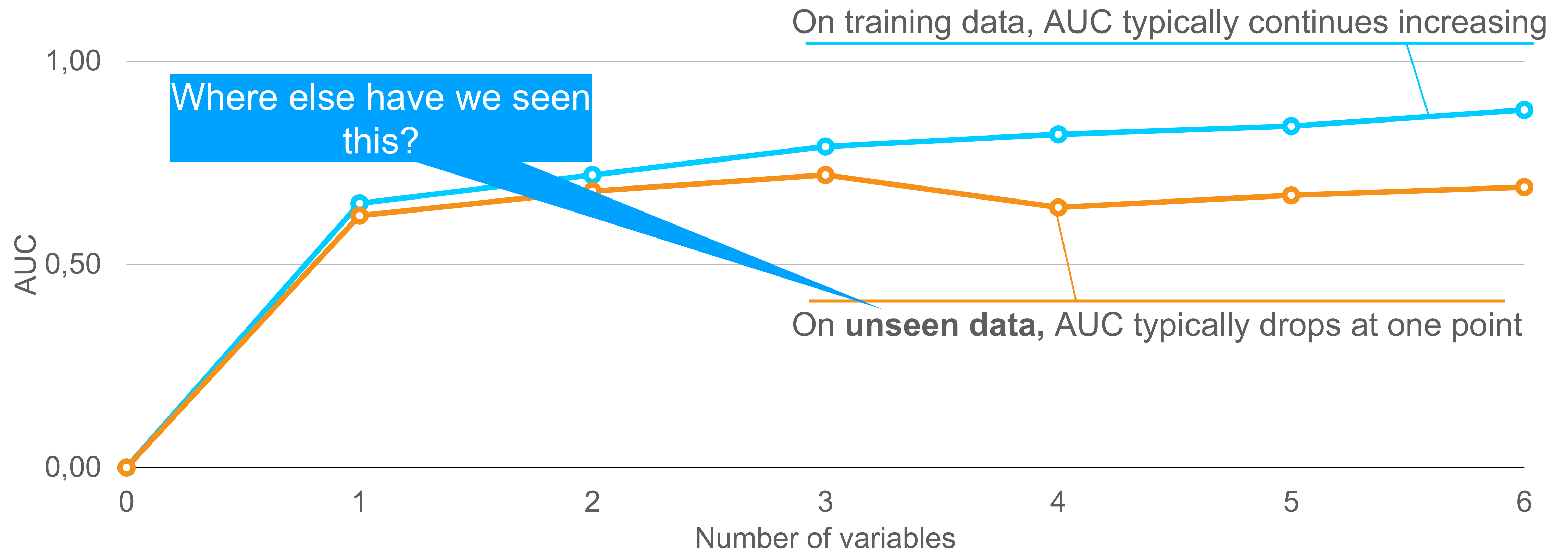| Variable 2<br>Variable 4 |

| Variable 2<br>Variable 4<br>Variable 6 |

| Variable 2<br>Variable 4<br>Variable 6<br>Variable 3 |

| Variable 2<br>Variable 4<br>Variable 6<br>Variable 3<br>Variable 1 |

| Variable 2<br>Variable 4<br>Variable 6<br>Variable 3<br>Variable 1<br>Variable 5 |

# Stepwise forward variable selection

# Stepwise forward variable selection



On training data, AUC typically continues increasing

On unseen data, AUC typically drops at one point

AUC

Number of variables

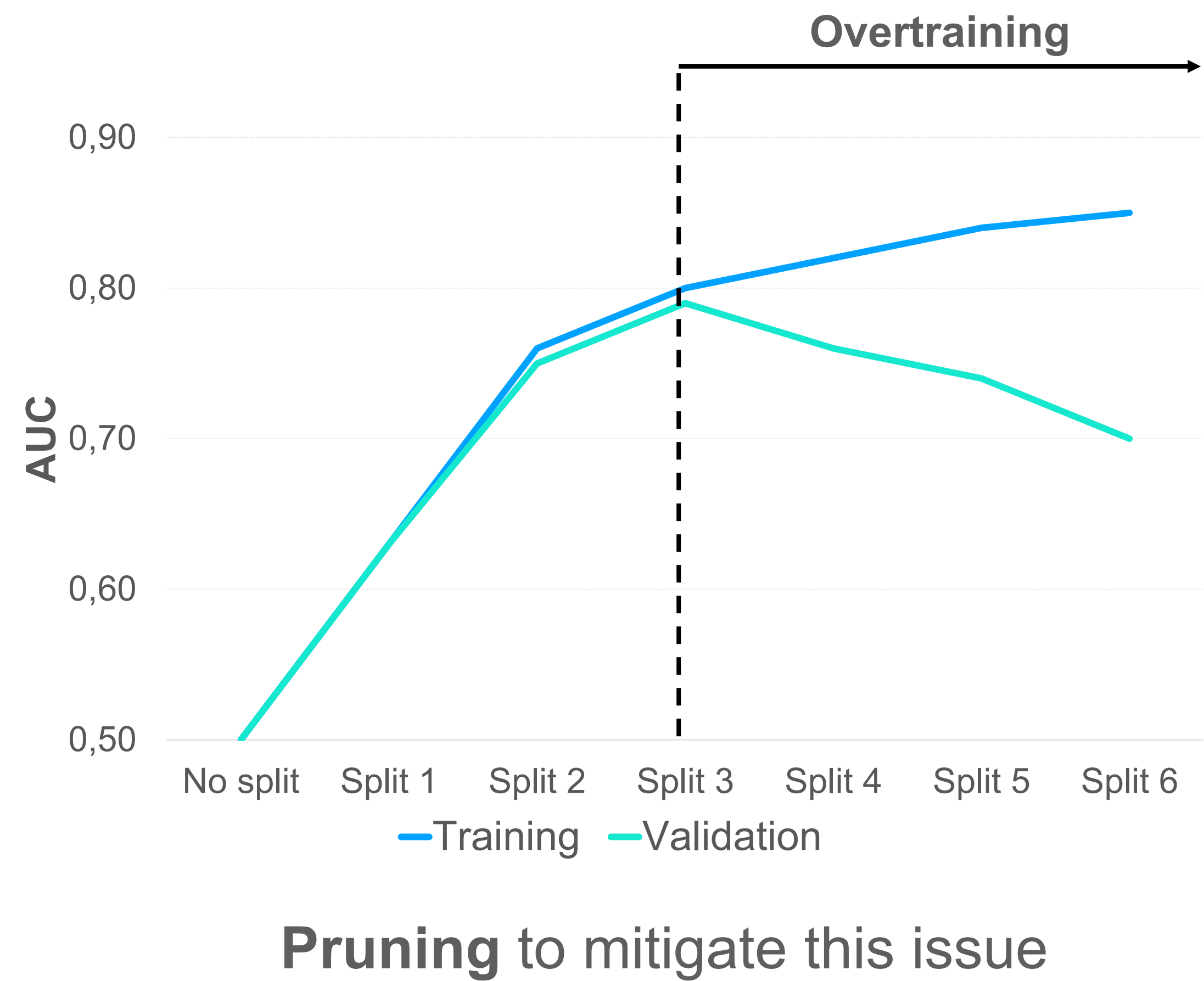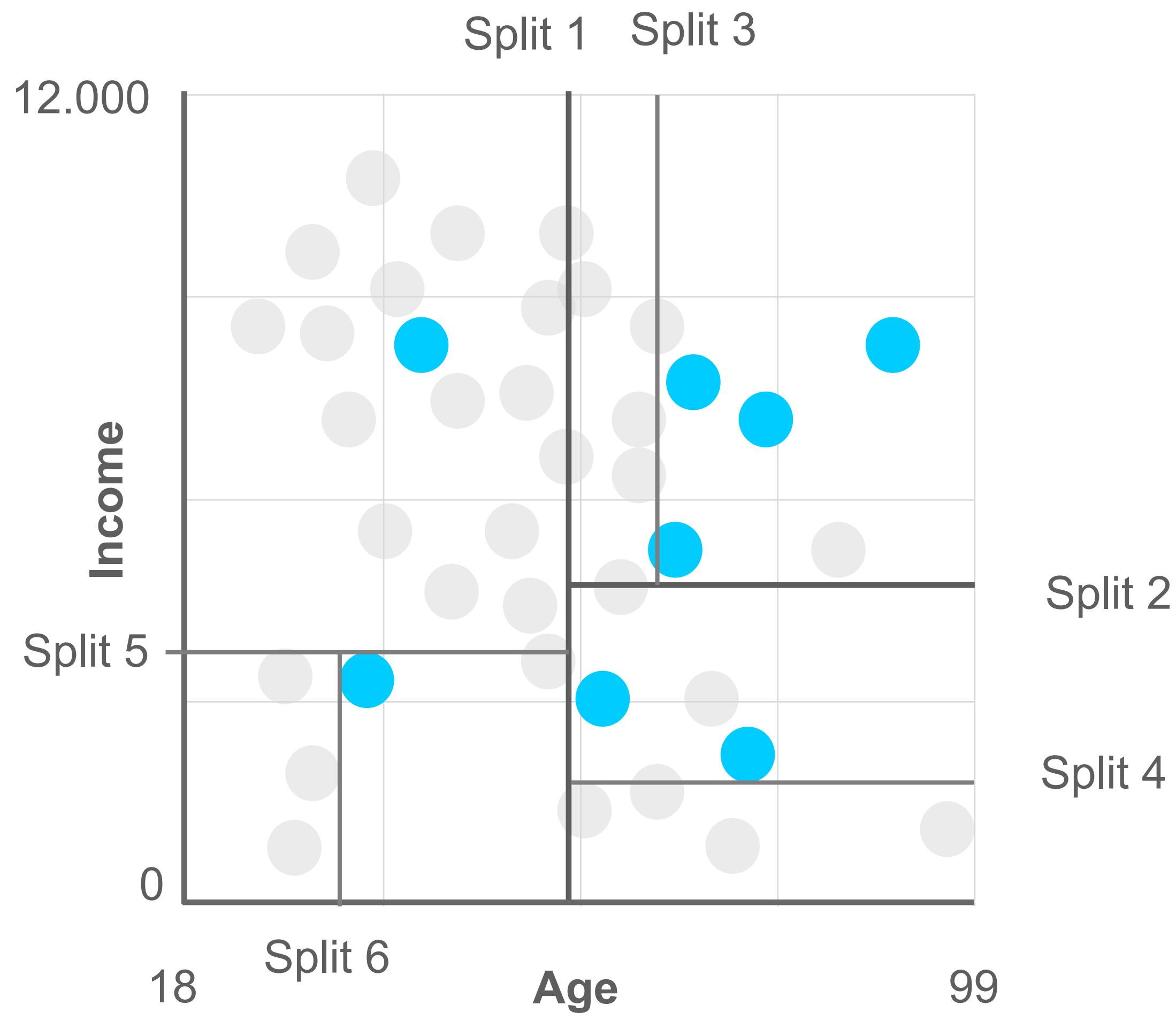# Univariate selection vs. Stepwise selection

- Will the outcome of selected features be the same?

# Importance of unseen data



On training data, AUC typically continues increasing

Where else have we seen this?

On **unseen data,** AUC typically drops at one point

AUC

0,00
0,50
1,00

0    1    2    3    4    5    6

Number of variables

# Intermezzo: let's get back to the decision tree

# Intermezzo: Decision tree



**Age**

< 55          ≥ 55

**9%**          **Income**

< 5.000          ≥ 5.000

**25%**          **44%**

**Maximum depth**

**Parameters**

- Optimized to minimize the impurity measurement
- Splitting decision at each branch (feature + value)

They are learned from the training data

**Hyperparameters**

- Configuration settings that control the behavior of the machine learning algorithm
- Maximum depth, Minimum samples per leaf,…

They are NOT learned from the training data
They are set prior to model training

**Control Model Complexity**
Hyperparameters can help manage the complexity of a model to **avoid overfitting or underfitting**

:

# Intermezzo: Regulating parameters

**https://scikit-learn.org/1.5/modules/generated/sklearn.svm.LinearSVC.html**

## LinearSVC

```
class sklearn.svm.LinearSVC(penalty='l2', loss='squared_hinge', *, dual='auto',
tol=0.0001, C=1.0, multi_class='ovr', fit_intercept=True, intercept_scaling=1,
class_weight=None, verbose=0, random_state=None, max_iter=1000)
```
[source]

**penalty : {'l1', 'l2'}, default='l2'**

Specifies the norm used in the penalization. The 'l2' penalty is the standard used in SVC. The 'l1' leads to `coef_` vectors that are sparse.

**C : float, default=1.0**

Regularization parameter. The strength of the regularization is inversely proportional to C. Must be strictly positive. For an intuitive visualization of the effects of scaling the regularization parameter C, see Scaling the regularization parameter for SVCs.

# Intermezzo: Regulating parameters

Regularization parameters play a crucial role in controlling the complexity of models and preventing overfitting across different algorithms:

**Linear Regression**: Regularization through Ridge, Lasso, or Elastic Net (parameters: λ\lambdaλ, λ1\lambda_1λ1, λ2\lambda_2λ2).

**Decision Trees**: Complexity controlled by parameters like max_depth, min_samples_split, min_samples_leaf, and max_features.

**Logistic Regression**: Regularization through Ridge, Lasso, or Elastic

By adjusting these regulating parameters, you can optimize your models for better performance and generalization on unseen data.